

Sémantique des déterminants dans un cadre richement typé

Christian Retoré

IRIT & Université de Bordeaux

`Christian.Retore@irit.fr`

7 février 2013

Résumé : La variation du sens des mots en contexte nous a conduit à enrichir le système de types utilisés dans notre analyse syntaxico-sémantique du français basé sur les grammaires catégorielles et la sémantique de Montague (ou la lambda-DRT). L’avantage majeur d’une telle sémantique profonde est de représenter le sens par des formules logiques aisément exploitables, par exemple par un moteur d’inférence. Déterminants et quantificateurs jouent un rôle fondamental dans la construction de ces formules. Mais dans notre système de types complexes, leurs termes sémantiques usuels ne fonctionnent pas. Nous proposons une solution inspirée des opérateurs epsilon et tau de Hilbert, sorte de fonctions de choix et d’éléments génériques. Cela unifie le traitement des différents types de déterminants et de quantificateurs ainsi que le liage dynamique des pronoms. Surtout, cette modélisation totalement calculable s’intègre parfaitement dans l’analyseur à large échelle du français Grail, tant en théorie qu’en pratique.

Mots-clefs : Analyse sémantique automatique ; Sémantique formelle ; Compositionnalité ;

Determiners in type theoretical semantics

Abstract : The variation of word meaning according to the context leads us to enrich the type system of our syntactical and semantic analyser of French based on categorial grammars and Montague semantics (or lambda-DRT). The main advantage of a deep semantic analyse is too represent meaning by logical formulae that can be easily used e.g. for inferences. Determiners and quantifiers play a fundamental role in the construction of those formulae. But in our rich type system the usual semantic terms do not work. We propose a solution inspired by the tau and epsilon operators of Hilbert, kinds of generic elements and choice functions. This approach unifies the treatment of the different determiners and quantifiers as well as the dynamic binding of pronouns. Above all, this fully computational view fits in well within the wide coverage parser Grail, both from a theoretical and a practical viewpoint.

Keywords : Automated semantic analysis ; Formal Semantics ; Compositional Semantics ;

1 Présentation

Dans le cadre du traitement automatique des langues, on entend plus souvent parler de sémantique distributionnelle, de vecteurs de mots et de fréquences que de sémantique formelle ou compositionnelle. Certes, les approches quantitatives sont plus aisées à mettre en oeuvre et fournissent des outils efficaces mais elles ne répondent pas aux mêmes questions. Les approches quantitatives sont fort utiles en recherche d'information et en classification car elles permettent de dire *de quoi parle* une phrase, une page web, un texte. En revanche elles ne disent pas ce qu'affirme le texte analysé, *qui fait quoi*. Une phrase peut très bien nier quelque chose, et ainsi causer une erreur à un système de recherche d'information (cf. exemple 1). Il faut aussi savoir reconnaître les pronoms, pour répondre à des questions comme *Geach était-il l'élève de Wittgenstein ?* à partir du web où on ne trouve que l'exemple 2. Dans une tâche de reconstruction d'itinéraires à partir de texte comme celui de l'exemple 3, il faut aussi analyser finement et complètement les parties pertinentes d'un texte, qui peuvent être trouvées par recherche d'information.¹

- (1) Mais vérification faite, ce n'était pas un ouragan qui était passé par là.
- (2) Bien qu'IL N'ait JAMAIS suivi l'enseignement académique de CE DERNIER, cependant IL EN éprouva fortement l'influence.
- (3) Le chemin pavé (...) serpente à travers fourrés de buis et de noisetiers. Puis, cinq minutes nous conduisent à un petit pont (...) qui nous porte sur la rive droite.

Ce montre que l'analyse sémantique complète et profonde d'une ou plusieurs phrases reste une tâche pertinente dans le traitement automatique des langues. Ce processus est utilement complété par des techniques statistiques, par exemple pour trouver les paragraphes pertinents ou pour définir des préférences contextuelles lorsque plusieurs sens sont possibles. Une fois construite avec ou sans l'aide de méthodes statistiques les représentations logiques du sens, j'affirme, au risque de froisser certains, que cela est bien supérieur aux graphes sémantiques, car les formules permettent de faire des inférences. Autrement dit, un moteur d'inférence utilisant les formules extraites du texte sera capable de dire si quelque chose découle de ce qui a été analysé.

Pour analyser logiquement une phrase, un ingrédient important est le traitement des déterminants. En effet, les déterminants indéfinis corresponde à une quantification existentielle (sauf le "un" générique, signifiant "tout"), tandis que les déterminants définis correspondent à la désignation d'un élément saillant du contexte. La quantification universelle est assez rare (et parfois juste exprimée par "les"), tandis que la quantification vague comme *beaucoup*, *peu*, *la plupart* est assez fréquente.

D'un point de vue pratique, ce travail se situe dans le cadre des grammaires catégorielles qui sont une approche de la syntaxe très orientée vers la sémantique compositionnelle. En effet, il est assez aisé de déduire de la structure syntaxique proposée par une grammaire catégorielle une représentation du sens sous forme logique. On observera que les deux systèmes produisant une analyse sémantique complète comme une formule logique (une DRS, en fait) sont basés sur les grammaires catégorielles : des CCG dans le cas de Boxer [5] et des MMCG dans le

1. Sauf mention contraire, nos exemples proviennent d'Internet.

ot	type sémantique u^* sémantique : λ -terme de type u^* x^v la variable ou la constante x est de type v
un	$(e \rightarrow t) \rightarrow ((e \rightarrow t) \rightarrow t)$ $\lambda P^{e \rightarrow t} \lambda Q^{e \rightarrow t} (\exists^{(e \rightarrow t) \rightarrow t} (\lambda x^e (\wedge^{t \rightarrow (t \rightarrow t)} (P\ x)(Q\ x))))$
club	$e \rightarrow t$ $\lambda x^e (\text{club}^{e \rightarrow t} x)$
a_battu	$e \rightarrow (e \rightarrow t)$ $\lambda y^e \lambda x^e ((\text{a_battu}^{e \rightarrow (e \rightarrow t)} x)y)$
Leeds	e Leeds

FIGURE 1 – Un lexique sémantique élémentaire

cas de [14, 15] que nous utilisons. Dans un cas comme dans l'autre la grammaire est acquise automatiquement sur corpus annoté. La première analyse l'anglais, la seconde le français. L'acquisition automatique de la grammaire produit un grand nombre de catégories par mot, et un minimum de traitement probabiliste est nécessaire pour ne considérer que les assignations les plus probables lors de l'analyse. Du point de vue sémantique, ces systèmes utilisent la correspondance entre syntaxe et sémantique telle qu'initée par Montague. On trouvera cette correspondance entre syntaxe et sémantique copieusement expliquée dans [18], mais nous allons la résumer brièvement avec un exemple, car c'est le point de départ de nos travaux.

Supposons que l'analyse syntaxique de "*un club a_battu Leeds*." produise "*(un (club)) (a_battu Leeds)*" expression dans laquelle la fonction est systématiquement écrite à gauche. Si les termes sémantiques sont ceux du lexique 1, alors en remplaçant les mots par les termes sémantiques associés on obtient un grand λ -terme, que l'on peut réduire :

$$\begin{aligned}
& \left((\lambda P^{e \rightarrow t} \lambda Q^{e \rightarrow t} (\exists^{(e \rightarrow t) \rightarrow t} (\lambda x^e (\wedge^{t \rightarrow (t \rightarrow t)} (P\ x)(Q\ x)))) (\lambda x^e (\text{club}^{e \rightarrow t} x)) \right) \\
& \quad \left((\lambda y^e \lambda x^e ((\text{a_battu}^{e \rightarrow (e \rightarrow t)} x)y)) \text{Leeds}^e \right) \\
& \quad \quad \downarrow \beta \\
& (\lambda Q^{e \rightarrow t} (\exists^{(e \rightarrow t) \rightarrow t} (\lambda x^e (\wedge^{t \rightarrow (t \rightarrow t)} (\text{club}^{e \rightarrow t} x)(Q\ x)))) \\
& \quad (\lambda x^e ((\text{a_battu}^{e \rightarrow (e \rightarrow t)} x)\text{Leeds}^e)) \\
& \quad \quad \downarrow \beta \\
& (\exists^{(e \rightarrow t) \rightarrow t} (\lambda x^e (\wedge^{t \rightarrow (t \rightarrow t)} x)((\text{a_battu}^{e \rightarrow (e \rightarrow t)} x)\text{Leeds}^e)))
\end{aligned}$$

Ce λ -terme de type **t** qui peut être appelé la forme logique de la phrase, est plus agréable sous un format standard : $\exists x : e (\text{club}(x) \wedge \text{a_battu}(x, \text{Leeds}))$.

On observera qu'il y a deux logiques à l'oeuvre : la première est le calcul propositionnel intuitionniste dont on n'utilise que les preuves ou λ -termes. Il servent à assembler des formules partielles. La seconde est une logique dont on n'utilise que les formules. Le terme de type **t** obtenu *in fine* est effectivement une formule logique, décrite en λ -calcul. Notre exemple contient un quantificateur, traité de manière classique : cela comporte quelques inconvénients, et surtout ce traitement est incompatible avec la prise en compte de la sémantique lexicale.

2 Déterminants et quantificateurs

Sans surprise, les déterminants considérés sont de deux sortes, définis et indéfinis — nous essaierons d’éviter les pluriels, qui posent d’autres problèmes. Bien évidemment, ce sont des constructions très fréquentes et fort importante dans la structure logique de la phrase. Nous allons essayer dans rendre compte, même si un travail de formalisation et d’automatisation comme celui-ci ne peut prétendre atteindre la finesse de [6]

Logiquement, les déterminants indéfinis correspondent à une quantification existentielle. La restriction du domaine de quantification, le nom commun qui suit le déterminant fait intervenir un second prédicat. Considérons quelques exemples :

- (4) J’ai senti **un** animal *me toucher le pied*.
- (5) **Un** parent d’élève de maternelle *vient chercher son enfant* en état d’ébriété, l’enseignant commet-il une faute en remettant l’enfant à ce parent ?
- (6) Aujourd’hui, je me suis réveillé en sursaut parce que j’ai senti **quelque chose** *me toucher le pied*. Il s’est avéré que c’était mon autre pied.
- (7) Précisez si **quelqu’un** *vient chercher l’enfant*.
- (8) Il y avait une panthère sortie de la cage. Elle était attachée. **L’animal** a sauté sur moi.
- (9) Un homme avait menacé la principale du collège de Monts où son fils était scolarisé. **Le parent d’élève** a été condamné hier.
- (10) Soudain, **un homme** est entré. **Il** a hurlé « Donne-moi la caisse ! ».
- (11) A la SPA si ont désire adopter **un animal** il faut donner 500F, et on a 24 Heures pour réfléchir si l’on désire **l’animal** ou non.

Les deux premiers exemples (4, 5) sont à mettre en parallèle avec les deux suivants (6, 7) qui correspondent eux-aussi à une quantification existentielle. Dans cette deuxième version, il n’y a que le prédicat principal, que nous avons choisi pour être le même, et il n’y a plus de restriction à une classe d’objets par un nom commun avec compléments (\bar{N} syntaxiquement ou $e \rightarrow t$ sémantiquement)

Les déterminants définis ont un rapport avec les déterminants indéfinis, qui souvent les introduisent, comme le montre les exemples (8,9) Les expressions se correspondent, et idéalement on aimerait que "*le X*" soit précédé de "*un X*", comme dans l’exemple 11. En fait, c’est plutôt rare, et les exemples en corpus sont plutôt comme 8 et 9 : l’antécédent de l’anaphore associative n’est pas celui qu’on espérerait pour un traitement automatique, quelques inférences sont nécessaires.

2.1 Traitement usuel et critique

L’analyse traditionnelle attribue à l’article indéfini un terme sémantique exprimant une quantification existentielle.

$$\begin{aligned} \text{un} : & \lambda P^{e \rightarrow t} \lambda Q^{e \rightarrow t} (\exists \lambda x^t. \&(P\ x)(Q\ x)) : (e \rightarrow t) \rightarrow (e \rightarrow t) \rightarrow t \\ \text{quelque chose} : & \exists : (e \rightarrow t) \rightarrow e \end{aligned}$$

Les articles définis sont traités différemment : les groupes nominaux qu'ils introduisent sont plutôt vus comme des anaphores, dites associatives, dont on cherche les référents. Une autre approche consiste à utiliser une fonction de choix, approche dont nous allons reparler. Cette modélisation classique en sémantique formelle ou dans les grammaires catégorielles pose divers problèmes.

Syntaxe et sémantique Les déterminants ainsi traités mettent à mal la correspondance entre syntaxe et sémantique. Cela oblige les grammaires catégorielles à s'éloigner des syntaxes usuelles et à avoir une catégorie pour chaque position syntaxique du GN quantifié.

- (12) elle écoutait une chanson de lassana hawa
- (13) SYNT. USUELLE : elle (écoutait (une (chanson (de lassana hawa))))
- (14) SEM. & CG : une (chanson de lassana hawa) (λx elle écoutait x)

Référence du GN quantifié Comme le fait remarquer [9], avant même que le prédicat principal soit énoncé, on peut se forger une interprétation du groupe nominal défini ou indéfini :

- (15) Ensuite, les élèves sont allés en salle info, pour réaliser un caryotype classé.
- (16) Ensuite, des élèves sont venus voir ce que l'on faisait.

Focus L'approche standard impose une symétrie entre le prédicat principal et la restriction à une classe d'objets que la langue ne fait pas.

- (17) Certains politiciens sont des menteurs car ce qui les intéresse (...)
- (18) * Certains menteurs sont des politiciens car ce qui les intéresse (...) ²

Définis et indéfinis Comme remarqué dans [7, 24, 25], l'unicité est loin d'être requise lorsque l'on utilise un déterminant indéfini. Un locuteur qui n'est pas au courant qu'il y en a trois, peut dire l'île du lac de Constance pour parler de celle qui voit. Effectivement on trouve :

- (19) Recueilli (...) par les moines de l'abbaye de Reichenau, sur **l'île du lac de Constance**,

De plus, "un" et "le" se rapprochent aussi car le contexte extra linguistique permet parfaitement d'utiliser l'article défini sans que le référent ait été introduit.

- (20) J'avais pris l'assurance 'automatiquement' avec le prêt immobilier lors de l'achat de *la maison*. ³

Bref, les déterminants "un" et se ressemblent, alors que les grammaires catégorielles et la sémantique formelle usuelle les opposent. Selon von Heusinger il s'agit d'une différence d'interprétation et non de forme logique : le "un" choisit un nouvel élément, tandis que le "le" choisit le plus saillant des référents possibles.

Pronoms de type E Une interprétation possible et très naturelle des pronoms due à [8] consiste interpréter le pronom par le terme sémantique de son antécédent, comme le ferait l'article défini, ce que font aussi [7, 24, 25]. Cela échappe au terme sémantique standard associé à "un" qu'on a vu ci-dessus.

2. Cet exemple est de nous, pour faire contraste avec le précédent.

3. Dans ce récit trouvé sur une FAQ il n'a jamais été question de "maison" auparavant

- (21) Soudain, un homme est entré.
Il / Cet homme / L'homme a hurlé « Donne-moi la caisse ! ».

3 Opérateurs de Hilbert, quantificateurs et déterminants

Les opérateurs de [10], surtout ι et ϵ ont été utilisés pour modéliser les déterminants et la quantification existentielle en particulier par von Steup [7, 24, 25], et plus récemment et sous une forme différente par [23].

Pour davantage de précision sur les opérateurs de Hilbert, on pourra consulter [22, 3]. [21] eut le premier l'idée d'introduire un *terme* — un individu — noté $\iota_x P(x)$ comme sémantique d'une description définie "*le P*" où "*P*" est une propriété, une formule à une variable libre. Que dénote ce terme ? Rien s'il n'existe pas un unique individu tel que $P(x)$, et sinon cet unique individu. Mais chacun sait que le quantificateur $\exists! x P(x)$ n'a pas de bonnes propriétés, notamment parce que son dual n'a rien de naturel : $(\forall x \neg P(x)) \vee (\exists y \exists z (y \neq z) \wedge P(y) \wedge P(z))$.

Hilbert s'est donc penché sur les termes génériques sans cette condition d'unicité. Il associe un terme $\epsilon_x F(x)$ à toute formule $F(x)$, et cette fois on a $F(\epsilon_x F(x)) \equiv \exists x F(x)$. Ce terme admet des règles de déduction relativement simples. De $F(t)$ on peut déduire $F(\epsilon_x F(x))$. Réciproquement si on a établi $F(x)$ sans x libre en hypothèse (1), on peut déduire $F(\epsilon_x \neg F(x))$: cette deuxième règle peut sembler compliquée, mais cela ne l'est pas tant que cela : (1) signifie qu'on a établi $\forall x P(x)$, soit $\neg(\exists x \neg P(x))$ c'est-à-dire $\neg(\neg P(\epsilon_x \neg P(x))) \equiv (\epsilon_x \neg F(x))$. Il y a aussi un terme dual $\tau_x P(x)$ avec les règles duales : de $F(x)$ sans x libre en hypothèse, on peut déduire $F(\tau_x F(x))$, et de $F(\tau_x F(x))$ on peut déduire $F(t)$ pour tout terme t . Bien sûr, vu cette dualité, un seul des deux opérateurs τ et ϵ suffit si on dispose de la négation. Pour les applications linguistiques on prend généralement ϵ car il y a bien plus de quantifications existentielles que de quantifications universelles dans la langue — c'est aussi pour cela que la DRT s'organise autour des quantifications existentielles.

L'idée est simplement de construire un terme générique associé au groupe nominal quantifié. Par exemple, pour "*un enfant sage*" on forme le terme $\epsilon_x.enfant(x) \& sage(x)$, et selon von Steup, pour "*l'enfant sage*" le terme est quasi identique : $\eta_x.enfant(x) \& sage(x)$. La différence entre η et ϵ n'est qu'une différence d'interprétation : ϵ choisit le plus présent en contexte tandis que η en choisit un nouveau. On remarque qu'on n'a pas parlé de typage, cependant, même si les auteurs mentionnés ne le font pas, ϵ est de type $(e \rightarrow t) \rightarrow e$: ϵ produit un individu (un terme) à partir d'une propriété.

4 Rappels sur le lexique génératif montagovien

Précédemment, nous avons proposé un lexique syntaxique et sémantique qui étend considérablement la sémantique de Montague pour rendre compte de l'adaptation du sens d'un mot au contexte et la coprédication [4]. Ce modèle s'est avéré pertinent pour des questions de sémantique lexicale ou compositionnelle : ambiguïté des déverbaux [19], le voyageur fictif [16], pluriels [17], termes génériques [20]). Ces travaux sont très comparables à ceux de [2] et de [12, 13], qui n'ont pas encore abordé la question des déterminants que nous étudions ici.

La question initiale est simple : comment modéliser les restrictions de sélection, comment rejeter les deux premiers exemples et accepter les suivants ?

- (22) * Une chaise aboie souvent.
- (23) * leur dix est bon
- (24) Mon chiot aboie souvent pour m'inciter à jouer avec.
- (25) Barcelone a battu Benfica 2-0.
- (26) Barcelone a choisi de structurer le réseau routier de manière à préserver un centre ville piétonnier.
- (27) Mon premier livre de cuisine ... Mon livre fétiche à cette époque !
- (28) Je l'ai retrouvé, il y a peu, chez ma maman [mon premier livre de cuisine]

L'idée est de typer les prédicats de sorte que les compositions sémantiquement impossible correspondent à un conflit de types. L'argument de aboyer doit être un chien, où tout au moins un animal, un nombre ne saurait appartenir à des personnes ni être bon, etc. L'impossibilité sémantique est matérialisée par l'application d'un prédicat $P^{\xi \rightarrow t}$ présupposant un argument de type ξ (par exemple animal) à un argument a^α de type α (par exemple, meuble) avec $\alpha \neq \xi$.

$$P^{\xi \rightarrow t} a^\alpha$$

On notera qu'il faut parfois relaxer ces contraintes : dans le contexte d'un match de rugby, "leur dix est bon" (exemple 23 ci-dessus) se trouve :

- (29) si leur dix est bon ils nous torchent ça c'est sûr

Il faut aussi prévoir que certaines coprédictions sont heureuses et d'autres moins :

- (30) ?? Barcelone a battu Benfica 2-0 et a choisi de structurer le réseau routier de manière à préserver un centre ville piétonnier.
- (31) Mon livre fétiche à cette époque, je l'ai retrouvé, il y a peu, chez ma maman.

Pour traiter tous ces phénomènes, nous avons proposé un lexique sémantique catégoriel où chaque mot se voit associer un λ -terme principal, qui ressemble beaucoup à celui de la sémantique de Montague rappelée ci-dessus, ainsi que des λ -termes optionnels qui permettent de transformer un mot dans l'aspect souhaité, par exemple un numéro en joueur de rugby etc.

Afin de factoriser les opérations sur des termes de types différents, ou d'avoir des opérations sur des familles de types, nous nous sommes placés dans le λ -calcul du second ordre appelé système F — mais des théories des types plus faibles, comme utilisés par Luo, seraient également possibles. En revanche, notre système se distingue surtout par le caractère lexical et non ontologique des transformations lexicales : celles-ci sont déclenchées par les mots et non par le type des mots. Cela nous semble pleinement justifié par des exemples comme "*promotion*" et "*classe*" : les deux désignent des groupes d'élèves, mais seul "*classe*", peut désigner un lieu (la salle de classe), tandis que "*promotion*" ne le peut pas.

Afin de rendre compte des coprédictions possibles ou impossibles, dans ce système où les mots portent les transformations, nous avons considéré deux

sortes de transformations. Celles-ci sont déclarées dans le lexiques comme rigides ou flexibles :

transformation rigide (F) Lorsqu'il y'a plusieurs occurrences de l'argument chacune la transformation rigide impose de n'utiliser qu'elle : le même aspect de l'argument doit être utilisé dans chacune des occurrences.

transformation flexible (R) Lorsqu'il y'a plusieurs occurrences de l'argument, une transformation différente peut être utilisée pour chaque occurrence. Cela permet de coprédiquer sur des aspects différents.

Il faut tout de même dire brièvement comment sont fait les λ -termes du second ordre, d'autant que la quantification sur les types joue un rôle dans notre traitement des déterminants. Les types du second ordre sont définis comme suit :

- Types de base :
 - \mathbf{t} les valeurs de vérité, \mathbf{v} les événements,
 - des types constants en grand nombre correspondant aux *différentes sortes d'individus*,
 - des *variables de type*, notées par des lettres grecques (issues d'un ensemble dénombrable P)
- Lorsque T est un type et α une variable de type, qui peut ou non apparaître dans T , $\Lambda\alpha. T$ est un type (dit polymorphe).
- Lorsque T_1 et T_2 sont des types, $T_1 \rightarrow T_2$ est aussi un type.

Pour définir les termes, on se donne une infinité dénombrable de variables de chaque type, ainsi que, pour chaque type, des constantes en nombre fini (possiblement aucune) :

- Une variable de type T c'est-à-dire $x : T$ (ce qu'on écrit aussi x^T) est un *terme* de type T .
- Une constante de type T c'est-à-dire $c : T$ (ce qu'on écrit aussi c^T) est un *terme* de type T .
- $(f \tau)$ est un terme de type U quant τ est de type T et f de type $T \rightarrow U$.
- $\lambda x^T. \tau$ est un terme de type $T \rightarrow U$ si x est une variable de type T , et τ un terme de type U .
- $\tau\{U\}$ est un terme de type $T[U/\alpha]$ quand $\tau : \Lambda\alpha. T$, et U est un type.
- $\Lambda\alpha. \tau$ est un terme de type $\Lambda\alpha. T$ quand α est une variable de type $\tau : T$ sans occurrence de α dans le type d'une variable libre.

Lorsque les constantes sont celles de la logique multisorte d'ordre supérieur (opérateurs $\& : \mathbf{t} \rightarrow \mathbf{t} \rightarrow \mathbf{t}, \forall : (e_i \rightarrow \mathbf{t}) \rightarrow \mathbf{t}, \dots$, et constantes du langage logique *regarde* : $(\mathbf{e}_a n i \rightarrow \mathbf{e} \rightarrow \mathbf{t})$ ce système est appelé $\Lambda\mathbf{Ty}_n$.

Les réductions pour λ et Λ sont définies de manière similaires.

- $(\Lambda\alpha. \tau)\{U\}$ se réduit en $\tau[U/\alpha]$ (rappelons que α et U sont des types).
- $(\lambda x. \tau)u$ se réduit en $\tau[u/x]$ (réduction habituelle).

La normalisation du système \mathbf{F} a une conséquence heureuse pour notre modèle sémantique : si les constantes (du λ -calcul) correspondent au langage L multisorte d'une logique d'ordre n (opérations logiques, prédicats, fonctions et constantes), tout terme normal de type \mathbf{t} correspond à une formule de L .

On précise ainsi l'organisation du système :

le λ -calcul du second ordre, le système \mathbf{F} sert à assembler les formules logiques partielles contenues dans le lexique (il remplace le λ -calcul simplement typé utilisé par Montague)

la logique d'ordre supérieur multisorte dans laquelle s'expriment les représentations sémantiques (elle remplace les formule de la logique d'ordre

word	principal λ -term	optional λ -terms	rigid/flexible
<i>ville</i>	$\widehat{T} : \mathbf{e} \rightarrow \mathbf{t}$	$Id_T : T \rightarrow T$ (F) $t_1 : T \rightarrow F$ (R) $t_2 : T \rightarrow P$ (F) $t_3 : T \rightarrow Pl$ (F)	
<i>Liverpool</i>	$liverpool^T$	$Id_T : T \rightarrow T$ (F) $t_1 : T \rightarrow F$ (R) $t_2 : T \rightarrow P$ (F) $t_3 : T \rightarrow Pl$ (F)	
<i>vaste</i>	$vaste : Pl \rightarrow \mathbf{t}$		
<i>a_voté</i>	$a_voté : P \rightarrow \mathbf{t}$		
<i>a_gagné</i>	$a_gagné : F \rightarrow \mathbf{t}$		

où les types de base sont définis comme suit T (ville), Pl (lieu), P (gens), F (club).

FIGURE 2 – Un exemple de lexique

supérieure utilisées par Montague, ou celle du premier ordre utilisé par réification : les nombreuses sortes sont les types de bases qui gère les restrictions de sélection).

Afin d'illustrer brièvement ce système et l'utilisation qui peut être faite de la quantification sur les types, donnons un exemple avec une coprédication qui fait intervenir une conjonction polymorphe. Etant donnés deux types α, β , deux prédicats $P^{\alpha \rightarrow \mathbf{t}}, Q^{\beta \rightarrow \mathbf{t}}$, sur des entités de sortes respectives α et β chaque fois qu'on aura un type ξ avec deux transformations de ξ dans α de ξ dans β , le système F contient un terme qui permet la coordination des propriétés P, Q de deux images d'une même entité de type ξ , et de le faire à chaque fois que l'on rencontrera pareille situation. "*Et*" polymorphe : $\&^\Pi =$

$$\Lambda \alpha \Lambda \beta \lambda P^{\alpha \rightarrow \mathbf{t}} \lambda Q^{\beta \rightarrow \mathbf{t}} \Lambda \xi \lambda x^\xi \lambda f^{\xi \rightarrow \alpha} \lambda g^{\xi \rightarrow \beta}. (\text{and}^{\mathbf{t} \rightarrow \mathbf{t} \rightarrow \mathbf{t}} (P (f x))(Q (g x)))$$

(32) *Liverpool* est *vaste* mais *a_voté*.

Cet exemple, s'analyse au moyen des deux transformations, celle d'une ville en un lieu et celle d'une ville en habitants. Aucune de ses deux transformations n'étant rigide on peut les utiliser simultanément. La conjonction polymorphe peut s'appliquer aux deux prédicats $(\&^\Pi (est_vaste)^{Pl \rightarrow \mathbf{t}} (a_voté)^{P \rightarrow \mathbf{t}})$ les variables de types sont alors instanciées par $\alpha := Pl$ et $\beta := P$ ce qui donne le terme $\&^\Pi \{Pl\} \{P\} (est_vaste)^{Pl \rightarrow \mathbf{t}} (a_voté)^{P \rightarrow \mathbf{t}}$ qui se réduit en $\Lambda \xi \lambda x^\xi \lambda f^{\xi \rightarrow \alpha} \lambda g^{\xi \rightarrow \beta} (\text{and}^{\mathbf{t} \rightarrow \mathbf{t} \rightarrow \mathbf{t}} (est_vaste (f x))(a_voté (g x)))$. La syntaxe nous conduit à appliquer ce terme à "*Liverpool*". ce qui impose l'instanciation $\xi := T$ et par réduction on obtient :

$\lambda f^{T \rightarrow Pl} \lambda g^{T \rightarrow P} (\text{and} (est_vaste (f Liverpool^T))(a_voté (g Liverpool^T)))$. Heureusement le lexique fournit deux λ -termes $t_2 : T \rightarrow P$ and $t_3 : T \rightarrow Pl$ qui peuvent être utilisés l'un et l'autre, puisqu'aucun des deux n'est rigide. Ainsi on obtient comme espéré

$$(\text{and} (est_vaste Pl \rightarrow \mathbf{t} (t_3^{T \rightarrow Pl} Liverpool^T))(a_voté^{Pl \rightarrow \mathbf{t}} (t_2^{T \rightarrow P} Liverpool^T)))$$

La même situation avec *a_voté* et *a_gagné* serait impossible car la transformation d'une ville en club est rigide.

5 Des termes typés pour les déterminants

5.1 Prédicats et types

Afin de traiter des exemples mentionnés il faut décider quels sont les types des prédicats puisqu'un opérateur de Hilbert se combine avec un prédicat pour donner un terme. Usuellement, un prédicat a pour type $\mathbf{e} \rightarrow \mathbf{t}$, mais en présence des innombrables types pour les entités autres que \mathbf{e} , que faire ? Faut-il autoriser des prédicats à avoir un domaine autre que \mathbf{e} ? Un prédicat comme le nom commun "*chat*" est-il une propriété des animaux ou une propriété de toutes les entités, propriété qui serait fausse en dehors des animaux ? Peu importe : un prédicat défini sur un type différent de \mathbf{e} , $P^{\alpha \rightarrow \mathbf{t}}$ s'étend en un $\overline{P}^{\mathbf{e} \rightarrow \mathbf{t}}$ sans difficulté, il suffit justement de dire qu'il est faux en dehors de α . Réciproquement, un prédicat comme *chat* défini sur α (par exemple animaux) peut se restreindre sur tout sous type de α . Evidemment, un prédicat comme *chat* restreint à un sous ensemble strict de l'ensemble où il est vrai (par exemple *siamois*) et ensuite étendu à \mathbf{e} puis restreint aux animaux ne redonnera pas le prédicat initial, car l'extension est définie uniformément sur tous les types et les prédicats comme étant fausse à l'extérieur du domaine considéré. Ainsi, lorsque β ne contient pas tout les $x : \alpha$ satisfaisant P on a $(\overline{P^{\alpha \rightarrow \mathbf{t}}|_{\beta}})|_{\alpha} \neq P$. On notera que c'est une question d'interprétation : dans le calcul des représentations sémantiques, seules l'extension et la restriction sont disponibles, pas leur interprétation dans un modèle.

On peut aussi se demander si un type définit un prédicat ? Si "*animal*" est un type, y a-t-il un prédicat "*être un animal*" ? Et si oui, quel est son domaine ? Etant donné un type α il est difficile de dire quel type β contenant α est un bon candidat pour le domaine du prédicat *être de type* α : aussi prendrons nous pour prédicat associé au type α le prédicat $\hat{\alpha}$ de type $\mathbf{e} \rightarrow \mathbf{t}$

Une question très délicate est celle des types de base.

- On pourrait n'avoir que \mathbf{e} et \mathbf{v} , mais cela ne permet pas la prise en compte de la sémantique lexicale.
- On pourrait considérer que tout formule à une variable libre définit un type. Cela en fait beaucoup, d'autant plus que comme les types permettent de définir des formules, il y a une sorte de circularité : comme on l'a vu, les types définissent naturellement des formules, et si celles-ci servent à définir des types, il n'est pas sûr que le système soit bien fondé.
- [2] propose d'utiliser un nombre de type de base relativement petit avec des types comme "*objet physique, contenu informationnel, humain,...*" correspondant aux restrictions de sélection que l'on rencontre dans la langue.
- Luo propose d'utiliser tous les noms communs.[13]
- Nous n'avons pas d'avis tranché sur la question, mais nous faisons remarquer qu'il faut sans doute aussi des types pour les propositions et pour les verbes d'action, car on quantifie aussi sur ce type d'objet :

(33) Elle voudrait qu'il croit en tout ce qu'elle lui dit.

(34) Il a fait tout ce qu'il a pu et il n'a même pas voulu être payé.

5.2 Des termes typés pour les déterminants

Nous proposons que les déterminants indéfinis soit modélisés par une constante de type $\Lambda\alpha. (\alpha \rightarrow \mathbf{t}) \rightarrow \alpha$. Nous voyons donc l'article indéfini comme un ϵ polymorphe, qui s'applique à un prédicat, et rend un objet du type auquel se prédicat s'applique. Considérons l'exemple suivant, très simple et inventé, afin d'illustrer notre propos :

(35) Un chat dort (sous ta voiture).

(36) $\text{un} : \Lambda\alpha. (\alpha \rightarrow \mathbf{t}) \rightarrow \alpha$

Supposons que "*chat*" soit un prédicat qui s'applique au type "*animal*". La variable de type α va s'instancier en "*animal*", le résultat sera "*un chat*" de type "*animal*" le prédicat "*dort (sous ta voiture)*" pouvant s'appliquer à un "*animal*", on obtiendra "*dort(un chat) .t*". C'est plutôt satisfaisant, mais rien ne dit que "*un chat*" ait la propriété d'être un chat ! Rien ne permet, dans la sémantique de "*un*" de produire cela. En revanche, lorsqu'on a décrit l'opérateur ϵ il a été dit que $P(\epsilon_x.P(x)) \equiv \exists x P(x)$. Comme lorsqu'on dit "*un*" chat dans le sens d'un chat particulier (et non d'un chat générique) celui-ci existe : nous ajoutons donc la présupposition $\text{chat}(\text{un chat})$ (c'est-à-dire $\text{chat}(\epsilon_x.\text{chat}(x))$). On notera que *un chat* étant de type "*animal*" le prédicat "*chat*" peut effectivement s'y appliquer. En revanche, il ne faut pas introduire systématiquement la présupposition $F(\epsilon_x.F(x))$ sans que "*un F*" ait été prononcé, sinon cela reviendrait à affirmer que toute propriété est satisfaite par au moins un individu.

Supposons maintenant que "*chat*" soit un type et non une propriété. Bien évidemment, on peut transformer le type "*chat*" en la propriété correspondante $\widehat{\text{chat}}$ comme expliqué au paragraphe ?? et procéder comme ci-dessus. Cependant on peut aussi utiliser une autre entrée pour le déterminant, à savoir $\text{un} : \Lambda\alpha. \alpha$. Ce type est celui du faux (\perp), mais on peut utiliser sans risque une constante de ce type, il ne s'agit pas d'un terme clos. Si on applique une telle constante au type "*chat*" obtient alors "*un chat*" de type chat sans même avoir à ajouter une présupposition, et comme le fait pertinemment remarquer [2] une déclaration de type $x : T$ est bien une forme de présupposition, car il est quasi impossible de la nier. Il n'y a aucune difficulté à appliquer le prédicat "*dort*" à ce chat, puisque les inclusions ontologiques font partie des transformations du lexique.

On peut traiter de la même manière les articles définis, comme le fait von Heusinger : seul le calcul de la référence sera différent. Tandis que l'article indéfini requiert un nouvel élément, l'article défini choisit au contraire un élément déjà présent en contexte. L'approche permet aussi de traiter les pronoms de type E de Evans. Le fait que les termes génériques soient typés n'y change rien. Pour interpréter les anaphores comme le "*il*" de l'exemple 10 ou 21 il suffit de recopier le terme sémantique associé à l'antécédent de "*il*".

Finalement, et c'est une très bonne nouvelle, la quantification universelle qui à notre connaissance n'a guère été étudiée dans ce cadre est extraordinairement simple : il suffit d'utiliser le terme générique construit avec τ , c'est-à-dire $\tau_x F(x)$. Ce dernier est bien plus facile à interpréter que $\epsilon_x.P(x)$: il s'agit de l'élément générique des démonstrations mathématiques, un objet qui, par rapport à F n'a pas de propriété : ainsi, s'il a la propriété F , tout le monde l'a. Il est donc interprété par un élément imaginaire, un générique universel.

6 Implémentation

Il n’y a pas de constructions spécifiques à ajouter à l’analyseur syntaxique et sémantique du français Grail pour le traitement que nous proposons des déterminants et des quantificateurs. L’extension pour prendre en compte notre nouveau style de lexique a déjà été testée, du moins sur de petits lexiques sémantiques que l’on a bien voulu saisir. En effet, la grammaire a été acquise sur corpus, mais à l’heure actuelle nul ne sait comment acquérir automatiquement les lexiques sophistiqués que nous utilisons. [15, 14]

Du point de vue syntaxique, les déterminants et quantificateurs ont une catégorie plus simple, et n’en n’ont qu’une : il n’est nul besoin de considérer autant de catégories syntaxiques qu’il y a de positions syntaxique possibles pour le déterminant : $gn \setminus n$ suffit. De ce point de vue, ce travail ressemble au lien que nous avons tissés entre grammaires génératives et grammaire catégorielle, et ce n’est pas un hasard, puisque c’était pour l’interprétation sémantique. [1]

Dans la pratique, plutôt que du λ -calcul nous utilisons de la λ -DRT pour calculer les représentations sémantiques. Cela change peu de choses, ce n’est pas plus difficile, mais nous préférons ne pas avoir à présenter la DRT et la λ -DRT. Cette variante permet à l’analyse sémantique de mieux suivre la structure discursive. Cela permet aussi de voir la correspondance dont parle [25] entre éléments génériques (l’approche présentée ici) et le liage dynamique des variables existentielles utilisé par la DRT.

7 Conclusion

Ce travail ouvre des questions en sémantique automatique et en logique.

Nous souhaitons explorer la portée des quantificateurs dans ce type d’approche, ainsi que le lien avec le calcul des prédicats dynamiques. Il semble que les opérateurs de Hilbert permettent de rapprocher ces deux visions assez différentes.

Nous avons évité de parler des pluriels qui sont bien sûr en rapport. Un premier travail a été effectué par [17] avec des opérateurs qui viennent formaliser des idées anciennes sur les pluriels. Néanmoins, le lien avec les déterminants et la quantification n’a pas été étudié. C’est sans doute une question de sémantique fort intéressante.

Une question pratique importante pour étendre l’analyseur sémantique est celle des types de base. Quels sont-ils ? Peut-être n’y a-t-il pas de réponse en général, car ils dépendent des restrictions de sélection dont on souhaite rendre compte, et du type d’informations attendues. Par exemple, dans le cadre de la reconstruction d’itinéraires entreprise à partir d’un corpus de récits de voyages du XIXe, des catégories spatiales et temporelles assez naturelles se dégagent. [11]

Notre travail pose également des questions d’acquisition, d’une part des types de base, mais aussi des λ -termes sémantiques. Sur la question particulière des déterminants il n’y a pas grand chose à dire : on les connaît, on connaît leur type leur sémantique associée, mais l’analyseur ne peut fonctionner que si les autres mots, noms, verbes, adjectifs ont des types : il faut donc connaître les types de base et les termes typés associés aux mots pour que l’analyseur fonctionne avec des lexiques plus gros que les tests effectués.

L'interprétation des termes génériques reste mystérieuse et pose des questions logiques amusantes. Il semble que le calcul avec ϵ ait été doté d'une sémantique, mais le livre n'est pas facile à trouver ni à lire. Les modèles sont inévitablement complexes puisque la logique avec ϵ permet d'exprimer des formules qui font pas partie des logiques usuelles, à rapprocher des quantificateurs branchants de Henkin. Par ailleurs, l'interprétation plus intuitive de von Heusinger, dépendante du contexte discursif, n'est peut-être pas pleinement satisfaisante non plus et on peut se demander si l'équivalence avec la quantification existentielle usuelle est toujours valide. [22, 3]

Toujours concernant la logique, mais aussi l'architecture du modèle informatique, l'interaction entre les types utilisés pour la composition du sens et les prédicats de la logique mutlisorte utilisé pour représenter le sens est très intrigante.

Références

- [1] Maxime Amblard, Alain Lecomte, and Christian Retoré. Categorical minimalist grammars : From generative grammar to logical form. *Linguistic Analysis*, 36(1–4) :273–306, 2010.
- [2] Nicholas Asher. *Lexical Meaning in context – a web of words*. Cambridge University press, 2011.
- [3] Jeremy Avigad and Richard Zach. The epsilon calculus. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Center for the Study of Language and Information, 2008.
- [4] Christian Bassac, Bruno Mery, and Christian Retoré. Towards a Type-Theoretical Account of Lexical Semantics. *Journal of Logic Language and Information*, 19(2) :229–245, April 2010. <http://hal.inria.fr/inria-00408308/>.
- [5] Johan Bos. Wide-coverage semantic analysis with boxer. In Johan Bos and Rodolfo Delmonte, editors, *Semantics in Text Processing. STEP 2008 Conference Proceedings*, Research in Computational Semantics, pages 277–286. College Publications, 2008.
- [6] Francis Corblin, Ileana Comorovski, Brenda Laca, and Claire Beyssade. Generalized quantifiers, dynamic semantics, and french determiners. In Francis Corblin and Henriëtte de Swart, editors, *Handbook of French Semantic*, chapter 1, pages 3–22. CSLI Publications, 2004.
- [7] Urs Egli and Klaus von Heusinger. The epsilon operator and E-type pronouns. In Urs Egli, Peter E. Pause, Christoph Schwarze, Arnim von Stechow, and Götz Wienold, editors, *Lexical Knowledge in the Organization of Language*, pages 121–141. Benjamins, 1995.
- [8] Gareth Evans. Pronouns, quantifiers, and relative clauses (i). *Canadian Journal of Philosophy*, 7(3) :467–536, 1977.
- [9] Peter Thomas Geach. *Reference and generality : an examination of some medieval and modern theories*. Contemporary philosophy. Cornell University Press, 1962.
- [10] David Hilbert and Paul Bernays. *Grundlagen der Mathematik. Bd. 2*. Springer, 1939. Traduction française de F. Gaillard, E. Guillaume et M. Guillaume, L'Harmattan, 2001.

- [11] Anaïs Lefeuvre, Richard Moot, Christian Retoré, and Noémie-Fleur Sandillon-Rezer. Traitement automatique sur corpus de récits de voyages pyrénéens : Une analyse syntaxique, sémantique et temporelle. In *Traitement Automatique du Langage Naturel, TALN'2012*, 2012.
- [12] Zhaohui Luo. Contextual analysis of word meanings in type-theoretical semantics. In Sylvain Pogodalla and Jean-Philippe Prost, editors, *LACL*, volume 6736 of *LNCS*, pages 159–174. Springer, 2011.
- [13] Zhaohui Luo. Common nouns as types. In Denis Béchet and Alexander Ja. Dikovsky, editors, *LACL*, volume 7351 of *Lecture Notes in Computer Science*, pages 173–185. Springer, 2012.
- [14] Richard Moot. Semi-automated extraction of a wide-coverage type-logical grammar for French. In *Proceedings of Traitement Automatique des Langues Naturelles (TALN)*, Montreal, 2010.
- [15] Richard Moot. Wide-coverage French syntax and semantics using Grail. In *Proceedings of Traitement Automatique des Langues Naturelles (TALN)*, Montreal, 2010.
- [16] Richard Moot, Laurent Prévot, and Christian Retoré. Un calcul de termes typés pour la pragmatique lexicale — chemins et voyageurs fictifs dans un corpus de récits de voyages. In *Traitement Automatique du Langage Naturel, TALN 2011*, pages 161–166, Montpellier, France, June 2011.
- [17] Richard Moot and Christian Retoré. Second order lambda calculus for meaning assembly : on the logical syntax of plurals. In *Coconat*, Tilburg, Pays-Bas, December 2011.
- [18] Richard Moot and Christian Retoré. *The logic of categorial grammars : a deductive account of natural language syntax and semantics*, volume 6850 of *LNCS*. Springer, 2012. <http://www.springer.com/computer/theoretical+computer+science/book/978-3-642-31554-1>.
- [19] Livy-Maria Real-Coelho and Christian Retoré. A generative montagovian lexicon for polysemous deverbal nouns. In *4th World Congress and School on Universal Logic – Workshop on Logic and linguistics.*, Rio de Janeiro, April 2013.
- [20] Christian Retoré. Variable types for meaning assembly : a logical syntax for generic noun phrases introduced by "most". *Recherches Linguistiques de Vincennes*, 41 :83–102, 2012. <http://hal.archives-ouvertes.fr/hal-00677312>.
- [21] Bertrand Russell. On denoting. *Mind*, 56(14) :479–493, 1905.
- [22] Barry Hartley Slater. Epsilon calculi. *The Internet Encyclopedia of Philosophy*, 2005.
- [23] Mark Steedman. *Taking Scope : The Natural Semantics of Quantifiers*. MIT Press, 2012.
- [24] Klaus von Heusinger. Definite descriptions and choice functions. In S. Akama, editor, *Logic, Language and Computation*, pages 61–91. Kluwer, 1997.
- [25] Klaus von Heusinger. Choice functions and the anaphoric semantics of definite nps. *Research on Language and Computation*, 2 :309–329, 2004.